

# The K giant stars from the LAMOST survey data I: identification, metallicity, and distance

Chao Liu<sup>1</sup>, Li-Cai Deng<sup>1</sup>, Jeffrey L. Carlin<sup>2</sup>, Martin C. Smith<sup>3</sup>, Jing Li<sup>1,3</sup>, Heidi Jo Newberg<sup>2</sup>, Shuang Gao<sup>1</sup>, Fan Yang<sup>1</sup>, Xiang-Xiang Xue<sup>4</sup>, Yan Xu<sup>1</sup>, Yue-Yang Zhang<sup>1</sup>, Yu Xin<sup>1</sup>, Ge Jin<sup>5</sup>

liuchao@nao.cas.cn

## ABSTRACT

We present a support vector machine classifier to identify the K giant stars from the LAMOST survey directly using their spectral line features. The completeness of the identification is about 75% for tests based on LAMOST stellar parameters. The contamination in the identified K giant sample is lower than 2.5%. Applying the classification method to about 2 million LAMOST spectra observed during the pilot survey and the first year survey, we select 298,036 K giant candidates. The metallicities of the sample are also estimated with uncertainty of  $0.13 \sim 0.29$  dex based on the equivalent widths of  $\text{Mg}_b$  and iron lines. A Bayesian method is then developed to estimate the posterior probability of the distance for the K giant stars, based on the estimated metallicity and 2MASS photometry. The synthetic isochrone-based distance estimates have been calibrated using 7 globular clusters with a wide range of metallicities. The uncertainty of the estimated distance modulus at  $K = 11$  mag, which is the median brightness of the K giant sample, is about 0.6 mag, corresponding to  $\sim 30\%$  in distance. As a scientific verification case, the trailing arm of the Sagittarius stream is clearly identified with the selected K giant sample. Moreover, at about 80 kpc from the Sun, we use our K giant stars to confirm a detection of stream members near the apo-center of the trailing tail. These rediscoveries of the features of the Sagittarius stream illustrate the potential of the LAMOST survey for detecting substructures in the halo of the Milky Way.

*Subject headings:* stars: K giants—stars: abundances—stars: distance—Galaxy: halo—Galaxy: structure

## 1. Introduction

The LAMOST (the Large sky Area Multi-Object fiber Spectroscopic Telescope; also known

as Guo Shou Jing Telescope) project has carried out a pilot survey between October 2011 and June 2012 and obtained more than 700,000 spectra (Cui et al. 2012; Zhao et al. 2012; Liu et al. 2013). The regular survey has operated since 2012 September and has already obtained about 2 million spectra as of June 2013. These spectra has been released as the DR1 catalog. A large fraction of them are K giant stars, which is of great interest in the studies of the Milky Way and in particular for the Galactic halo.

K giant stars are luminous and thus allow us to probe the Galaxy far beyond the solar neighborhood. Typically, the absolute magnitude of K giant stars is between  $M_r = 2$  and  $-2$  mag. Given

<sup>1</sup>Key Laboratory of Optical Astronomy, National Astronomical Observatories, Chinese Academy of Sciences, Datun Road 20A, Beijing 100012, China

<sup>2</sup>Department of Physics, Applied Physics and Astronomy, Rensselaer Polytechnic Institute, 110 8th Street, Troy, NY 12180, USA

<sup>3</sup>Shanghai Astronomical Observatory, Chinese Academy of Sciences, 80 Nandan Road, Shanghai 200030, China

<sup>4</sup>Max Planck Institute for Astronomy, Königstuhl 17, Heidelberg D-69117, Germany

<sup>5</sup>University of Science and Technology of China, Hefei 230026, China

the limiting magnitude of  $r = 17.8$  mag, the maximum distance at which LAMOST can detect K giant stars is  $\sim 90$  kpc from the Sun.

Xue et al. (2014) found distances to more than 4000 K giant stars in SDSS/SEGUE survey with accuracy of  $\sim 12\%$ . Because the limiting magnitude of SEGUE spectra is  $g = 20.2$  mag, the maximum distance is much larger than that of LAMOST. Although the K giant stars observed by LAMOST cannot compete with SDSS/SEGUE in terms of distance probed, the total number of this type of spectra in LAMOST will be two orders of magnitude larger than those in SDSS due to the huge number of spectra observed in the LAMOST survey and the fact that a larger fraction of bright stars are giants.

With such a huge dataset, we should be able to address many interesting and important questions on our galaxy, e.g., the total mass of the Milky Way, the shape of the dark matter halo, the kinematic substructures in the stellar spheroid, the mass distribution of the Galactic disk, the chemo-dynamical features and the evolution history of the disk, etc. Specifically, it will significantly improve the observational evidence of the kinematic substructures in the stellar halo (Deng et al. 2012). In the 2MASS and SDSS surveys, many substructures have been discovered in the past decades, including the Sagittarius dwarf galaxy stream (Ibata et al. 2001; Newberg et al. 2002; Majewski et al. 2003), Monoceros ring (Newberg et al. 2002; Yanny et al. 2003), Orphan stream (Belokurov et al. 2006), Virgo overdensity (Newberg et al. 2002; Vivas & Zinn 2006; Newberg et al. 2007), Triangulum-Andromeda overdensity (Majewski et al. 2004b; Rocha-Pinto et al. 2004), Hercules-Aquila cloud (Belokurov et al. 2007), Cetus polar stellar stream (Newberg, Yanny, & Willman 2009), Pisces stellar stream (Bonaca, Geha, & Kallivayalil 2012; Martin et al. 2013), and many other cold and weak streams (e.g., Gillmair & Dionatos 2006). Although some of these substructures are prominent in photometric catalogs, a tiny fraction of their member stars have spectroscopic observations. The identification of the K giant members in known tidal streams will be crucial to constrain the orbits of the tidal streams and to constrain the merging history of their progenitors (e.g., Law, Johnston, & Majewski 2005; Law & Majewski 2010). Moreover, it can also be

used to measure the total mass of the Milky Way (e.g., Koposov, Rix, & Hogg 2010).

In addition, the K giant stars can be used to discover new substructures which are otherwise not possible by any previous approaches. According to the  $\Lambda$ CDM cosmology, the halo of a Milky Way-like galaxy should contain hundreds of subhaloes, in which dwarf galaxies may be embedded. Current observations only find a few tidal streams and dozens of satellite dwarf galaxies around the Galaxy. This is the so called *missing satellite problem*, which challenges all current theories (Klypin et al. 1999; Koposov et al. 2008). Because some of the merging dwarf galaxies form tidal streams during accretion, the search for new tidal substructure is one important way to address this discrepancy. Ultimately, the study of tidal streams will allow a better understanding on the formation history and evolution of a galaxy.

In order to make optimal use of the K giant sample, a clean and relatively complete K giant catalog with distance estimates is highly desirable. In principle, K giant stars can be identified from measurements of stellar parameters, e.g. effective temperature ( $T_{\text{eff}}$ ) and surface gravity ( $\log g$ ). However, estimations of the stellar parameters can only be reliably achieved for the *good* spectra, i.e., the high signal-to-noise ratio data or the well flux calibrated data. As a result, a majority of K giant spectra with moderate or low signal-to-noise ratio, which are usually located at further distances, are missing from the parameter table. In order to reach a larger detection volume and hence maximize the scientific value of the LAMOST survey data, we take an alternative approach of identifying K giant stars directly from spectra instead of stellar parameters. In addition to the identification of the K giant stars, the metallicity is necessary for the distance estimation. Consequently, one of the aims of the current work is to develop a more robust and reliable metallicity estimation method, especially for those spectra with low or moderate S/N. Finally, the distance of the K giant samples is determined with a robust statistical method.

A brief introduction to the data used in this work is given in section 2. In section 3, a machine learning algorithm is developed for the identification of K giant stars from the stellar spectra of the LAMOST survey. The success of the classifica-

tion method is verified by using MILES and SDSS public data. The method is then applied to the LAMOST data and a complete K giant dataset is produced. In section 4, a thin-plate spline model (hereafter, LM2D) for [Fe/H] estimation based on the  $Mg_b$  and iron lines is introduced and applied to the identified K giant stars. Subsequently, the distance of the selected K giant stars with low extinction is estimated based on the isochrone comparison in section 5. The well-known kinematic substructure, the Sagittarius stream, is then identified from the K giant stars in section 6. Finally, a short conclusion of the current work is given in the last section.

## 2. LAMOST survey data

The LAMOST DR1 catalog contains more than 2 million spectra, including about 700,000 spectra observed during the pilot survey. In total, there are about 1.9 million stellar spectra available from the LAMOST survey.

The standard LAMOST pipeline (Luo et al. 2012, Luo et al. in preparation) converts 2D into 1D spectra, corrects the flat field, combines the blue and red parts of the spectra, calibrates the wavelength and subtracts the sky background. Figure 1 shows sample spectra for K giant stars processed by the standard pipeline.

For all spectra, the standard pipeline also provides the radial velocity based on the cross correlation with ELODIE library (Prugniel et al. 2007). However, only about half of the total data are high quality F, G, and K type spectra, for which stellar parameters are estimated using *Ulyss* (Wu et al. 2011) in the pipeline. *Ulyss* is a forward model method, which models the spectrum pixels as a linear combination of a set of non-linear components. Each non-linear component is a function of stellar parameters as well as radial velocity and it is defined in advance based on the ELODIE stellar library. The best fit stellar parameters and radial velocity of an observed stellar spectrum is determined iteratively by minimizing the  $\chi^2$  value between the observed spectrum and the model. Because the stellar parameters of only about 50% spectra with higher signal-to-noise ratio have been estimated, the selection function of the stars may be distorted and the sampling power of the survey may be weakened

## 3. K giant selection

### 3.1. Support vector machine classifier

Support vector machine (SVM) is a machine learning algorithm which is suited for classification (Cortes & Vapnik 1995) and broadly used in astronomy (e.g., Bailer-Jones et al. 2008; Liu et al. 2012; Saglia et al. 2012; Bailer-Jones et al. 2013). As a supervised algorithm, it needs a set of known data to train the SVM model first. A subset of the training data are selected as the support vectors during the training phase. The support vectors define the linear boundary of classes in a high dimensional inner product space. When the training process is done, the SVM model is ready for prediction; any data input to the model will be marked as a certain class depending on the region to which the input data are projected.

### 3.2. Data preprocessing

The full spectrum of a star does not carry useful information in every pixel. Some parts of the spectrum that contain strong line features play more important role in classification and parametrization. Therefore, we use the equivalent width ( $EW$ ) of such strong lines to identify the K giant stars.

First, we select some of the Balmer lines,  $H_\alpha$ ,  $H_\beta$ ,  $H_\gamma$ , and  $H_\delta$ , as the indicators of temperature. However, for stars with  $T_{\text{eff}}$  lower than  $\sim 5000$  K, the Balmer lines are very weak or even invisible, so the TiO feature near  $6150\text{\AA}$  is used as temperature indicator in this case. Second, in order to distinguish the giant from dwarf stars, we need to use the lines that are sensitive to surface gravity. The magnesium lines around  $5180\text{\AA}$  (including  $Mg_1$ ,  $Mg_2$ , and  $Mg_b$ ) are good tracers for this purpose. Some other important features, e.g. CN, G band, are also included. Finally, we use the  $EW$ s of the ten selected lines to characterize the whole spectrum. We adopt the Lick indexes (Worthey et al. 1994) to measure the  $EW$ s using the following equation:

$$EW = \int (1 - \frac{f_{\text{line}}(\lambda)}{f_{\text{cont}}(\lambda)}) d\lambda, \quad (1)$$

where  $f_{\text{cont}}(\lambda)$  and  $f_{\text{line}}(\lambda)$  are the fluxes of the continuum and the spectral line, respectively, both of which are functions of the wavelength  $\lambda$ .

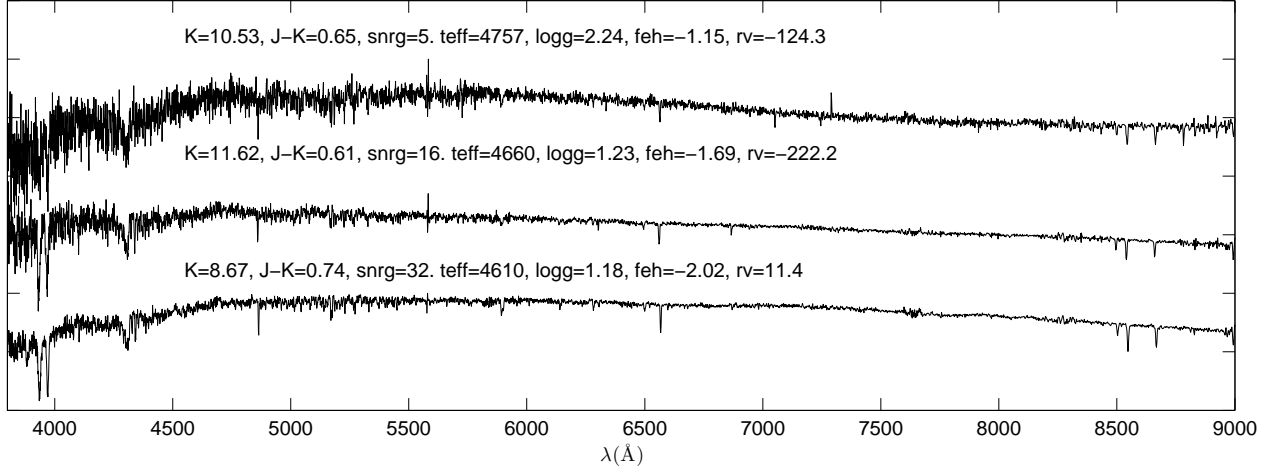


Fig. 1.— Sample K giant spectra selected from LAMOST survey. The parameters,  $K$  magnitude,  $J - K$ , signal-to-noise ratio in  $g$  band (snrg), effective temperature (teff), surface gravity (logg), metallicity (feh), and radial velocity (rv), are marked above the corresponding spectrum.

Figure 2 shows the difference between the K giant and non-K giant stars of the MILES library (Sánchez-Blázquez et al. 2006) in the  $EW$ s of  $Mg_b$ ,  $H_\beta$ , and  $TiO$ . We can essentially separate the K giant stars from  $EW_{H_\beta}$  vs.  $EW_{Mg_b}$  with some local overlapping.  $TiO$  can help to distinguish the giant from the dwarf stars in some of the overlapped region, particularly at  $4 < EW_{Mg_b} < 5$ .

### 3.3. Training of the SVM classifier

By training with a set of data with known giant/non-giant separation, the parameters within SVM can be properly tuned to obtain the best model for the classification. For this purpose, the training data is defined using a common sample of the LAMOST pilot survey and SDSS DR9 (Ahn et al. 2012) with the following additional criteria: i) the signal-to-noise ratio for SDSS spectra are larger than 20 and  $T_{\text{eff}}$ ,  $\log g$ , and  $[Fe/H]$  is provided by SSPP pipeline; ii) the spectra is marked as STAR in the LAMOST pipeline; and iii) the signal-to-noise ratio in  $g$  band (denoted as  $S/N(g)$ ) measured from LAMOST spectra is larger than  $10^1$ . Totally 2,046 matched objects

are selected.

In this sample, the *true* K giant stars are defined as:  $\log g < 4$  when  $4600 < T_{\text{eff}} < 5600$  or  $\log g < 3.5$  when  $T_{\text{eff}} < 4600$  (see the green polygon shown in figure 3). It is noted that the SDSS SSPP pipeline does not fit M stars, i.e.,  $T_{\text{eff}} < 4000$ , hence the SVM classifier trained by SDSS stellar parameters is also not suited for M giant stars. There are 274 *true* K giant stars defined in the training dataset. The rest of the 2,046 stars in the sample are marked as non-K giant stars. The  $EW$ s of the spectral lines are measured from the corresponding LAMOST spectra of the training sample so that the signal to noise level, wavelength calibration, residual of sky subtraction etc. are comparable to the real dataset.

We use the *libsvm* package<sup>2</sup> in MATLAB to train the SVM classifier based on the training sample.

### 3.4. Validation and performance of the classifier

Two different samples are selected as the test data to validate the classification. First, we use

<sup>1</sup>The SDSS  $g$  band filter covers from 4000 to 5300 Å, containing most of the useful spectral features, e.g., a few Balmer lines, CH and CN lines, Mg triplet, lots of iron lines etc., for K giant identification as well as parametrization. There-

fore, we use the  $S/N$  in  $g$  band to quantify the quality of the spectra.

<sup>2</sup><http://www.csie.ntu.edu.tw/~cjlin/libsvm>



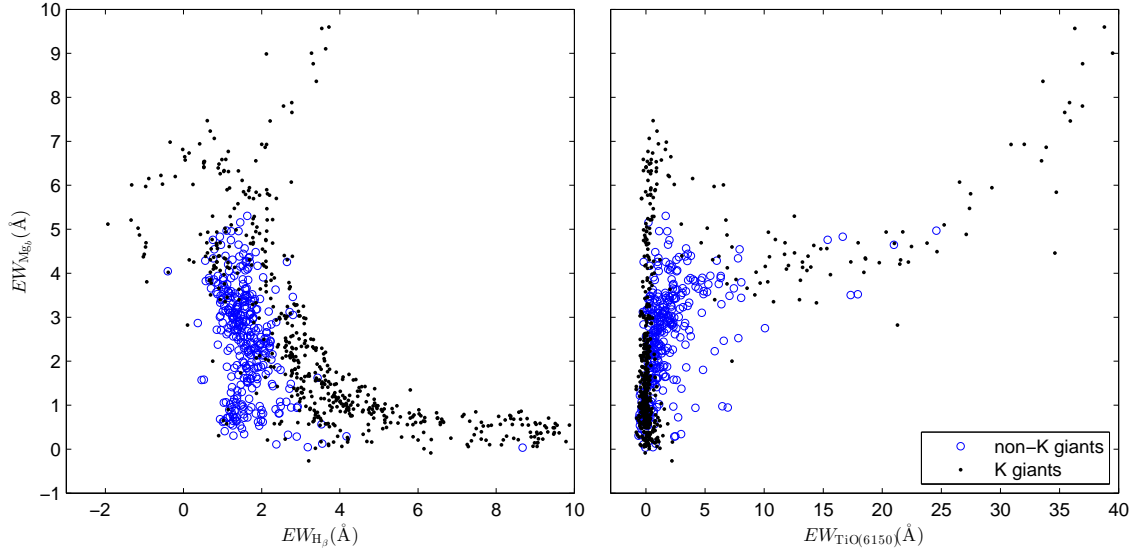


Fig. 2.— The figure shows the  $EW$ s of  $Mg_b$ ,  $H_\beta$ , and  $TiO$  ( $6150\text{\AA}$ ) for the MILES K giant stars (blue circles) and non-K giants (black dots). The left panel is  $EW_{H_\beta}$  vs.  $EW_{Mg_b}$  and the right panel is  $EW_{TiO(6150)}$  vs.  $EW_{Mg_b}$ .

the MILES library, which provides high S/N low-resolution spectra. Then we test the performance using the LAMOST spectra combined with the *true* K giant labels defined in SDSS stellar parameters. Table 1 summarizes the results of the performance tests.

The MILES library contains 985 low-resolution spectra in total. Note that stars with  $T_{\text{eff}} > 10000\text{K}$  are not suitable for the Lick indexes because their Balmer lines are too broad. After eliminating these, 919 stars are left as the test dataset. There are 350 stars defined as the *true* K giant stars with  $4000 < T_{\text{eff}} < 5600\text{K}$  and  $\log g < 4$  (see the green rectangle shown in figure 4). It is noted that the K giant selection criteria for SDSS and MILES data are different. This is because that in SDSS parameter plane (see figure 3) the  $\log g$  of the late-type main-sequence stars ( $T_{\text{eff}} < 5000\text{K}$ ) decreases for unknown reason. Hence, the criteria of the K giant selection based on SDSS parameters avoid these dwarf stars. For the MILES case, as shown in figure 4 the main-sequence is quite normal in low-temperature end. Therefore, we simply use a rectangle to select the *true* K giant stars. The different criteria may not lead to significant inconsistency, because the two *true* K

giant samples selected from the MILES data based on the two criteria are almost the same. The SVM classifier identifies 325 K giant stars, out of which 304 are *true* K giant stars. The *completeness* is defined as the ratio of the number of *identified* to that of the total *true* K giant samples, while the *contamination* as the fraction of the non-K giant stars in the *identified* K giant samples. Under these definitions, the completeness of the K giant stars identification for MILES library is 86.9% and the contamination is only 6.5%. The performance of the K giant selection is very good as shown in  $T_{\text{eff}}-\log g$  diagram in figure 4.

It is also noted that the MILES spectra have very high signal-to-noise ratio. In order to test the performance for lower quality data, such as the LAMOST spectra, artificial noise is added to the  $EW$ s to simulate the realistic situation. The classification algorithm is run 10 times with 10% random Gaussian noise<sup>3</sup> added on the  $EW$ s and obtain the values of the mean completeness of 84.1% and contamination of 6.7%; when the noise level goes up to 20% the mean completeness of 77.1%

<sup>3</sup>It means that the sigma of the Gaussian is 10% of the  $EW$  of a line. It is equivalent with  $S/N=10$ .

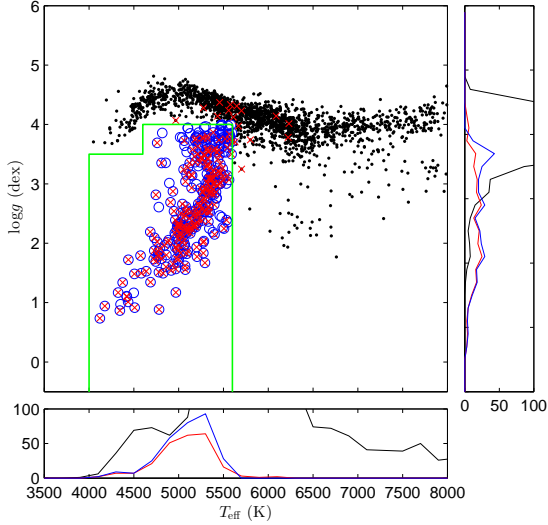


Fig. 3.— Distribution in  $T_{\text{eff}}$  vs.  $\log g$  for the LAMOST-SDSS cross-matched stars, for which  $T_{\text{eff}}$  and  $\log g$  are from the SDSS SSPP pipeline. The blue circles are the *true* K giant stars and the black dots the non-K giant stars. The red crosses indicate the *identified* K giant stars from the SVM classifier. The right and bottom panels show the histograms of  $\log g$  and  $T_{\text{eff}}$ , respectively. The black, blue, and red curves are the stellar counts for the full dataset, the *true*, and the *identified* K giant stars, respectively. The green rectangle shows the selection criteria of the *true* K giant stars.

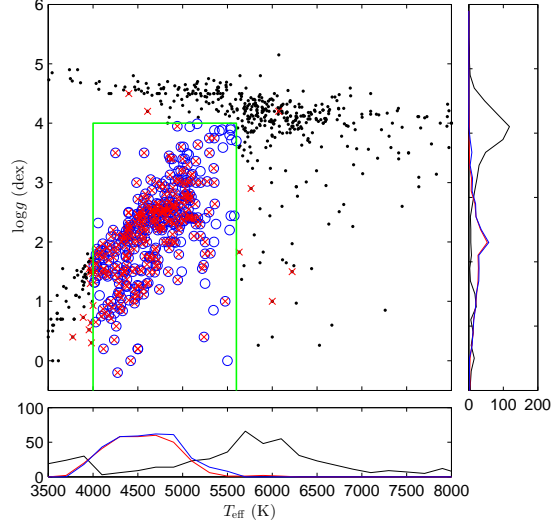


Fig. 4.— Distribution in  $T_{\text{eff}}$  vs.  $\log g$  for the MILES library. The symbols and the sideways histograms are the same as in Figure 4.

and contamination of 8.3% are obtained. Thus, we are convinced that the SVM classifier is robust even for low S/N spectra.

The second test is carried out using LAMOST regular survey data instead of MILES. We select 2,251 common stars in both LAMOST and SDSS DR9 catalog complying with the criteria in section 3.3. These test data are more representative than the MILES library, because they are observed and reduced in very similar situations to those for the training dataset. We define 302 *true* K giant stars in the test sample using the same definition in section 3.3. We then apply the SVM classifier to them and classify 236 K giant stars with 220 *true* K giant stars and 16 contaminates. In other words, the completeness of the classification is 72.8% and the contamination is 6.8%. Figure 3 shows the performance of the K giant selection in  $T_{\text{eff}}$ - $\log g$  diagram.

To investigate how the signal-to-noise ratio of the LAMOST spectra affect the classification, we separate the test data into two groups at  $S/N(g)=20$ . For those with  $S/N(g)<20$ , the completeness and the contamination are 72.1% and 9.6%, respectively, while the two values turn out to be 74.0% and 2.2% for those with  $S/N(g)>20$ . Even though the LAMOST spectra are more af-

affected by noise, the completeness of the test based on the LAMOST spectra with SDSS parameters is only  $\sim 10\%$  lower than that of the first test based on MILES. The result is very promising in the sense that even in the very tough case, i.e., low S/N spectra, more than 70% of the K giant stars can still be identified.

### 3.5. Application to the LAMOST survey

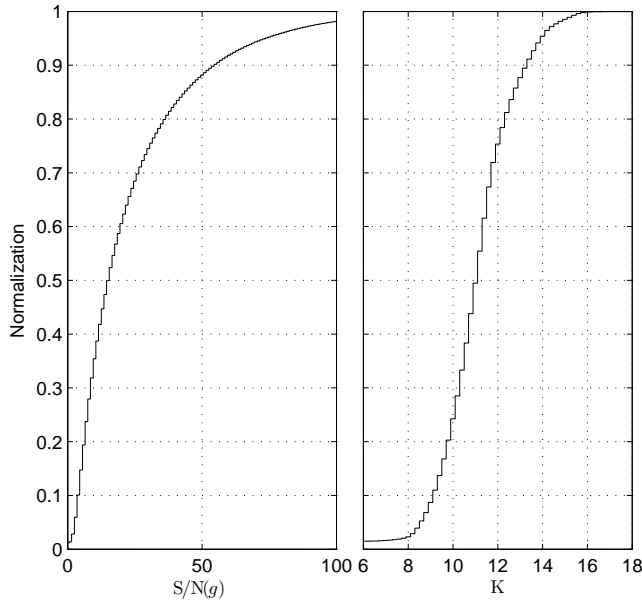


Fig. 5.— *Left panel:* The cumulative distribution of  $S/N(g)$  for the 298,036 identified K giant stars. *Right panel:* The cumulative distribution of 2MASS  $K$  magnitude for the identified K giant stars.

We apply the SVM classifier to all  $\sim 1.9$  million stellar spectra from the LAMOST DR1 survey data. The training process took about one day in a DELL workstation, while the process identification of K giant star in the whole dataset was very fast, taking about two hours with 12 parallel threads. In total, we obtain 298,036 identified K giant spectra of which 196,440 (119,813) have  $S/N(g) \geq 10$  ( $\geq 20$ ). The left panel of figure 5 shows the distribution of the  $S/N(g)$  for the identified K giant stars; and the right panel shows the distribution of 2MASS  $K$  magnitude for the same sample. The median value of  $S/N(g)$  is about 15

and the median  $K$  magnitude for the samples is at 11 mag.

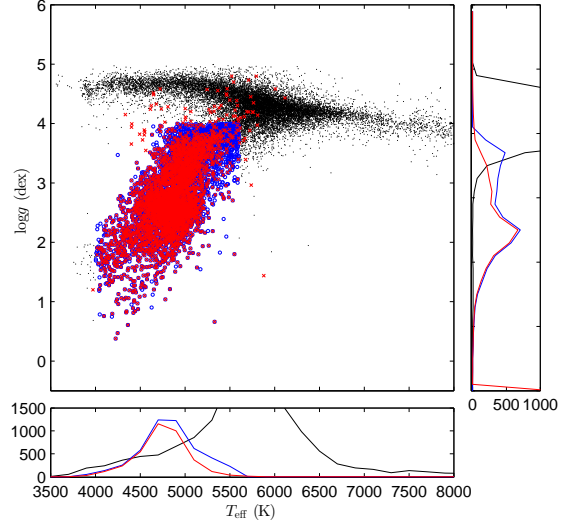


Fig. 6.— Distribution in  $T_{\text{eff}}$  vs.  $\log g$  of the LAMOST stars, for which  $T_{\text{eff}}$  and  $\log g$  are given by LAMOST pipeline. The symbols are the same as in Figure 4. In order to demonstrate the performance, only arbitrary one fiftieth spectra are drawn in the figure.

The LAMOST pipeline also provides stellar parameters for a selection of about 1 million high quality spectra. Among them, there are 238,332 *true* K giant spectra based on the definition given in section 3.3 but with stellar parameters provided by LAMOST pipeline. Applying the classification algorithm to the sample, the completeness and contamination becomes 74.5% and 2.4%, respectively, which is slightly better than the test with SDSS parameters (see Table 1). When dividing the data into two groups at  $S/N(g)=20$ , the completeness values of the low and high S/N data are 66.7% and 79.3%, contamination is 5.0% and 0.9%, respectively. Figure 6 shows the selected K giant stars in  $T_{\text{eff}}-\log g$  diagram.

It is clear that the classification algorithm is more successful for high quality data (high S/N); spectral quality of the survey is one of the keys to success.

Table 1: The performance of the SVM K giant classifier.

test data	completeness	contamination
MILES	86.9%	6.5%
MILES (10% noise) <sup>a</sup>	84.1%	6.7%
MILES (20% noise)	77.1%	8.3%
LAMOST+SSPP params. <sup>b</sup>	72.8%	6.8%
LAMOST+SSPP params. (S/N( <i>g</i> )<20)	72.1%	9.6%
LAMOST+SSPP params. (S/N( <i>g</i> )>20)	74.0%	2.2%
LAMOST+LAMOST params. <sup>c</sup>	74.5%	2.4%
LAMOST+LAMOST params. (S/N( <i>g</i> )<20)	66.7%	5.0%
LAMOST+LAMOST params. (S/N( <i>g</i> )>20)	79.3%	0.9%

<sup>a</sup>Add 10% Gaussian noise to the *EW*s of the MILES spectra.

<sup>b</sup>The test dataset contains the LAMOST spectra with SSPP derived stellar parameters.

<sup>c</sup>The test dataset contains the LAMOST spectra with LAMOST derived stellar parameters.

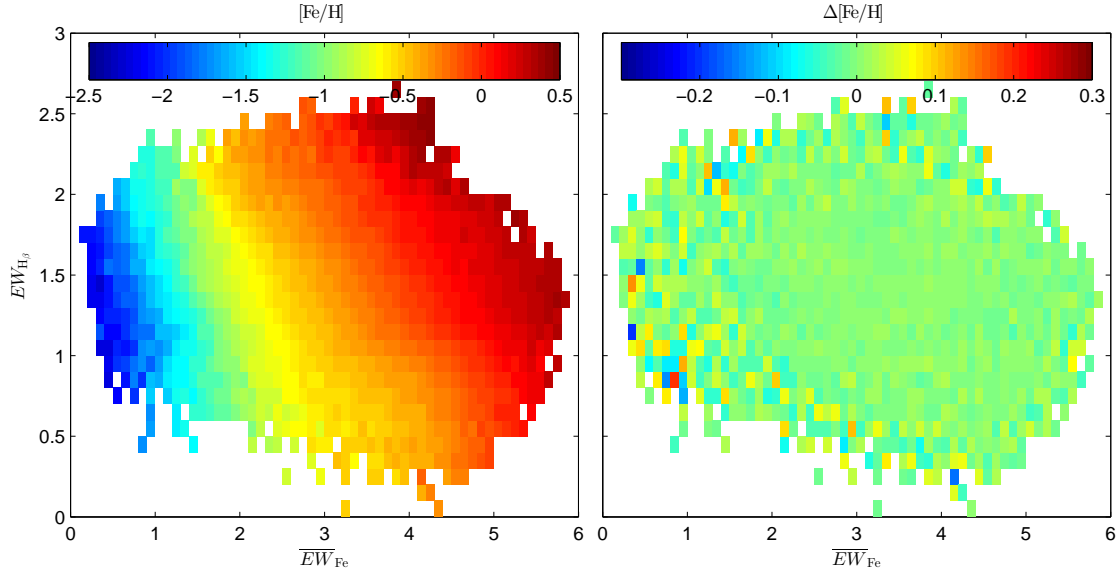


Fig. 7.— *Left panel*: The color codes the median LAMOST  $[\text{Fe}/\text{H}]$  in  $\overline{EW}_{\text{Fe}}$  vs.  $EW_{\text{H}\beta}$  plane for 114,900 K giant stars with LAMOST  $[\text{Fe}/\text{H}]$ . *Right panel*: The residual of the LM2D estimated  $[\text{Fe}/\text{H}]$ .

#### 4. Metallicity estimation

In general, the distance of a giant star is a function of magnitude, color index, and metallicity. Therefore, the metallicity of K giant stars has to be determined before any reliable distance can be derived.

The LAMOST pipeline has provided metallicity estimates for about 1 million high quality spectra. However, there are lots of K giant stars, which are mostly low S/N data, selected in section 3.5 with no estimated metallicity. In order to estimate the metallicity for all identified K giant stars including those with low S/N spectra, we use the LAMOST derived metallicity (denoted as  $[\text{Fe}/\text{H}]_{LM}$ ) as the training dataset and establish an estimator (hereafter named LM2D) to estimate  $[\text{Fe}/\text{H}]$  from the line features.

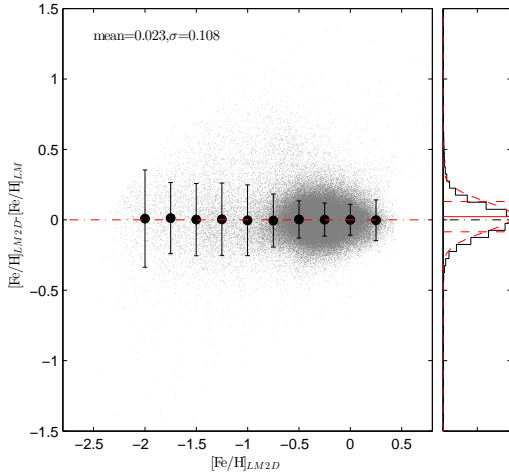


Fig. 8.— The residual of the  $[\text{Fe}/\text{H}]$  estimation from LM2D compared with the LAMOST  $[\text{Fe}/\text{H}]$ . The gray dots stand for the individual stars. The big black dots with error bars are the median residual at each  $[\text{Fe}/\text{H}]$  bin. The distribution of the residual is shown in the right side panel as the black line. The red dashed curve shows the best Gaussian fit of the distribution and the solid and dashed horizontal lines show the mean and  $1\sigma$  values, respectively.

In general, the metallicity of the K giant stars is a function of both the Fe lines and the effective temperature. Therefore, we simultaneously select  $EW_{H\beta}$ , which is correlated with  $T_{\text{eff}}$ , and

$\overline{EW}_{\text{Fe}}$ , which is the mean  $EW$ s of 9 Fe lines at 4383Å, 4531Å, 4668Å, 5015Å, 5270Å, 5335Å, 5406Å, 5709Å, and 5782Å, to estimate  $[\text{Fe}/\text{H}]$ .

The training dataset is selected with the following criteria: i) they are labeled as K giant stars based the LAMOST  $T_{\text{eff}}$  and  $\log g$  values according to section 3.3; ii) they are also identified as K giant stars in the SVM classifier; iii) they have  $S/N(g) > 20$ ; iv) they have  $[\text{Fe}/\text{H}]$  measured by LAMOST; and v)  $EW_{H\beta} > 0$ . There are 114,900 K giant stars that meet all of these conditions.

The median  $[\text{Fe}/\text{H}]$  of the sampled stars forms a curved surface in  $EW_{H\beta}$  vs.  $\overline{EW}_{\text{Fe}}$  plane (see the left panel of figure 7), which can be fitted with a 2-dimensional thin-plate spline function<sup>4</sup> depending on both  $EW_{H\beta}$  and  $\overline{EW}_{\text{Fe}}$ . The best fit thin-plate spline is then used to predict the metallicity of the K giant stars. The right panel of figure 7 shows the residuals of the spline fitting for the 114,900 K giant stars. For most of the area the residual is smaller than 0.1 dex. Figure 8 shows  $[\text{Fe}/\text{H}]_{LM2D} - [\text{Fe}/\text{H}]_{LM}$  as a function of  $[\text{Fe}/\text{H}]_{LM}$  for the training dataset. Although the dispersion of the residual of  $[\text{Fe}/\text{H}]$  becomes larger at low metallicity end, no significant systematics is found in figure 8. The overall dispersion of the differential  $[\text{Fe}/\text{H}]$  is 0.11 dex, being smaller for metal-rich stars but larger for metal-poor stars. The very weak line features in metal-poor spectra are responsible for the poor measurements of the equivalent widths and consequently account for the larger uncertainty in the metallicity estimation.

The comparisons of the derived metallicity  $[\text{Fe}/\text{H}]_{LM2D}$  with  $[\text{Fe}/\text{H}]_{LM}$  and SDSS metallicity,  $[\text{Fe}/\text{H}]_{SDSS}$ , will provide more tests of the performance of the LM2D method. We start by collecting the test data sample. The LAMOST dataset is cross-matched with SDSS DR9 and cleaned using the following criteria: 1) the S/N of the SDSS spectra is higher than 20; 2) The stellar parameters are provided by SDSS; 3) they are K giant stars according to criteria defined in section 3.3 using the SDSS parameters, and 4) they are also identified as K giant stars with the SVM classifier. A sample of 636 K giant stars with  $[\text{Fe}/\text{H}]_{LM2D}$  are resolved, of which 394 also have

<sup>4</sup>The thin-plate spline is a generalized spline used for two or more dimensional interpolation and smoothing.

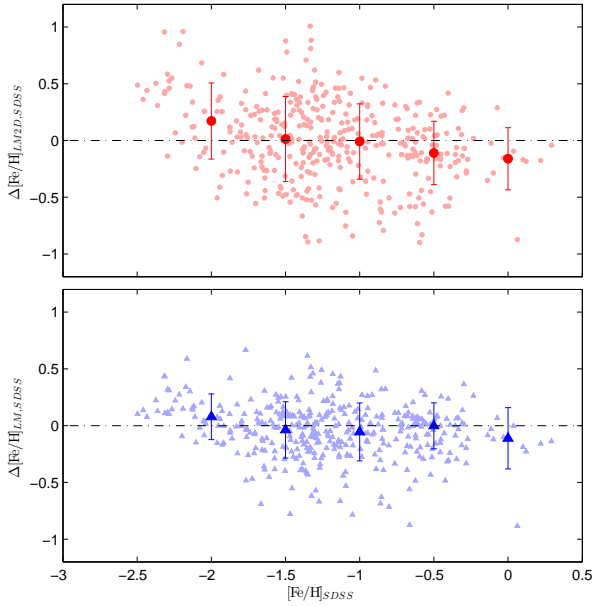


Fig. 9.— *Top panel:* The difference between  $[\text{Fe}/\text{H}]_{\text{LM2D}}$  and  $[\text{Fe}/\text{H}]_{\text{SDSS}}$  as a function of  $[\text{Fe}/\text{H}]_{\text{SDSS}}$  for 394 K giant stars with LM2D, LAMOST, and SDSS metallicities (pale red dots). The filled red circles with error bars show the medians and standard deviations at  $[\text{Fe}/\text{H}]=-2, -1.5, -1, -0.5,$  and  $0$ . *Bottom panel:* The difference between  $[\text{Fe}/\text{H}]_{\text{LM}}$  and  $[\text{Fe}/\text{H}]_{\text{SDSS}}$  as a function of  $[\text{Fe}/\text{H}]_{\text{SDSS}}$  (pale blue triangles). The blue triangles with error bars are the medians and standard deviations.

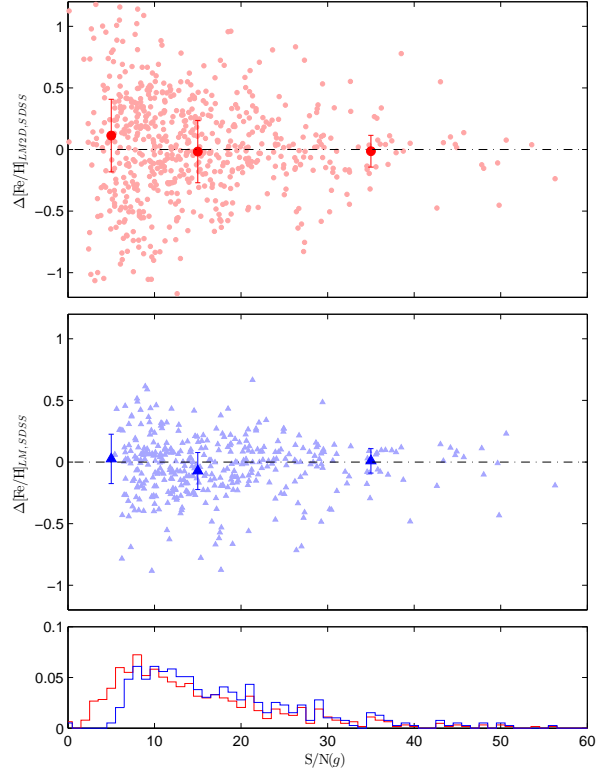


Fig. 10.— *Top panel:* The difference between  $[\text{Fe}/\text{H}]_{\text{LM2D}}$  and  $[\text{Fe}/\text{H}]_{\text{SDSS}}$  as a function of  $g$ -band signal-to-noise ratio,  $S/N(g)$ , in LAMOST spectra for 636 K giant stars with both LM2D and SDSS metallicities (pale red dots). The filled red circles with error bars show the medians and the median absolute deviations at  $S/N(g)=0-10, 10-20,$  and  $> 20$ . *Middle panel:* The difference between  $[\text{Fe}/\text{H}]_{\text{LM}}$  and  $[\text{Fe}/\text{H}]_{\text{SDSS}}$  as a function of  $S/N(g)$  for 394 K giant stars with both LAMOST and SDSS metallicities (pale blue triangles). The blue triangles with error bars are the median and standard deviation values. *Bottom panel:* The red (blue) line shows the histogram of the  $S/N$  for the K giant stars with LM2D (LAMOST) metallicity.

LAMOST  $[\text{Fe}/\text{H}]_{LM}$ . Notice that the SDSS data are dominated by metal-poor stars due to their fainter magnitude limit along with the fact that they are at higher latitude (see figure 9). Hence, the uncertainty of the LM2D metallicity is larger than the average level, as expected from figure 8.

Figure 9 shows the comparison between  $[\text{Fe}/\text{H}]_{LM2D}$  ( $[\text{Fe}/\text{H}]_{LM}$ ) and  $[\text{Fe}/\text{H}]_{SDSS}$  using the 394 stars with both  $[\text{Fe}/\text{H}]_{LM2D}$  and  $[\text{Fe}/\text{H}]_{LM}$ . In the top panel, it is indicated that the LM2D overestimates the  $[\text{Fe}/\text{H}]$  for stars with  $[\text{Fe}/\text{H}]_{SDSS} < -1$  by about 0.2 dex and underestimates the value by the similar level at  $[\text{Fe}/\text{H}]_{SDSS} = 0$ . As a comparison, the bottom panel shows that the LAMOST  $[\text{Fe}/\text{H}]$  has similar trend with smaller bias and dispersion.

Table 2: The performance of the  $[\text{Fe}/\text{H}]_{LM}$  and  $[\text{Fe}/\text{H}]_{LM2D}$  compared with SDSS metallicity.

S/N( <i>g</i> )	0-10	10-20	> 20
$\Delta[\text{Fe}/\text{H}]_{LM,SDSS}^a$	0.03	-0.07	0.01
$\sigma_{LM,SDSS}^b$	0.20	0.15	0.10
$\Delta[\text{Fe}/\text{H}]_{LM2D,SDSS}^c$	0.11	-0.02	-0.01
$\sigma_{LM2D,SDSS}^d$	0.29	0.25	0.13
$\Delta[\text{Fe}/\text{H}]_{LM2D,LM}^e$	0.00	0.04	-0.02
$\sigma_{LM2D,LM}^f$	0.28	0.20	0.10

<sup>a</sup>  $\Delta[\text{Fe}/\text{H}]_{LM,SDSS} = [\text{Fe}/\text{H}]_{LM} - [\text{Fe}/\text{H}]_{SDSS}$

<sup>b</sup> The median absolute deviation of  $[\text{Fe}/\text{H}]_{LM} - [\text{Fe}/\text{H}]_{SDSS}$ .

<sup>c</sup>  $\Delta[\text{Fe}/\text{H}]_{LM2D,SDSS} = [\text{Fe}/\text{H}]_{LM2D} - [\text{Fe}/\text{H}]_{SDSS}$ .

<sup>d</sup> The median absolute deviation of  $[\text{Fe}/\text{H}]_{LM2D} - [\text{Fe}/\text{H}]_{SDSS}$ .

<sup>e</sup>  $\Delta[\text{Fe}/\text{H}]_{LM2D,LM} = [\text{Fe}/\text{H}]_{LM2D} - [\text{Fe}/\text{H}]_{LM}$ .

<sup>f</sup> The median absolute deviation of  $[\text{Fe}/\text{H}]_{LM2D} - [\text{Fe}/\text{H}]_{LM}$ .

Figure 10 demonstrates the influence of the S/N of the LAMOST spectra in the LM2D metallicity estimation. The top panel presents the difference between  $[\text{Fe}/\text{H}]_{LM2D}$  and  $[\text{Fe}/\text{H}]_{SDSS}$  as a function of S/N(*g*). Not surprisingly, the dispersion of the difference of metallicity declines when the signal-to-noise ratio increases. The middle panel shows the same trend in LAMOST  $[\text{Fe}/\text{H}]_{LM}$ , although the dispersions are smaller than those in  $[\text{Fe}/\text{H}]_{LM2D}$ . As shown in the bottom panel, it is clear that LM2D can be applied to less preferable data, i.e., S/N(*g*) can be as low as 3, though the measured metallicity will be somewhat less accu-

rate.

Table 2 shows the offsets and dispersion of  $[\text{Fe}/\text{H}]_{LM}$  ( $[\text{Fe}/\text{H}]_{LM2D}$ ) compared with  $[\text{Fe}/\text{H}]_{SDSS}$  as a function of S/N(*g*). Compared with SDSS parameters, the uncertainty of LM2D method for K giant stars with S/N(*g*) < 10 is only about 0.1 dex larger than that of the LAMOST metallicity. And the uncertainty of LM2D is similar as that of LAMOST metallicity for the spectra with S/N(*g*) > 20.

## 5. The estimation of distance

The estimation of the distance for K giant stars is a non-trivial task, because the K giant stars are distributed along the red giant branch, therefore the luminosity is a steep function of effective temperature or color index.

Xue et al. (2014) estimated the distance with accuracy down to 0.2 mag in distance moduli (DM) using cluster-based fiducial isochrones. Although the clusters provide more realistic isochrones than synthetic data and hence match the observed data, the 4 clusters and one BaSTI (Pietrinferni et al. 2004) isochrone (at  $[\text{Fe}/\text{H}] = 0$ ) they used are only sparsely located at 5 discrete values in  $[\text{Fe}/\text{H}]$ . The estimated distance is very sensitive to the accuracy of the  $[\text{Fe}/\text{H}]$ , particularly for metal-poor stars. Thus, a denser grid in  $[\text{Fe}/\text{H}]$  may improve the accuracy of distance.

In this section, we estimate the distance of the halo K giant stars in an alternative way, i.e., using synthetic isochrones with calibration based on globular clusters. Because the disk K giant stars have much broader range of age, their distance estimation needs different method and beyond this work.

In principle, the synthetic isochrones can provide arbitrary grid in  $[\text{Fe}/\text{H}]$ . In our case, the mean uncertainty of the metallicity is  $\sim 0.3$  dex. Therefore, we use the synthetic library with  $\Delta Z = 0.0001$ , which corresponds to the  $\Delta[\text{Fe}/\text{H}] \sim 0.3$  dex at  $[\text{Fe}/\text{H}] = -2$  dex. To correct the discrepancy between the synthetic isochrones and the real observed objects, we calibrate the synthetic data using globular clusters.

### 5.1. Method

Because there is no unified accurate photometry covering a large range of  $r = 10\text{--}18$  mag,

which matches the LAMOST spectra, the photometry of the targets are composed of different catalogs. UCAC4 is used for the targeting brighter than  $r = 14$  mag, while SDSS and PanSTARRS are used for the targeting fainter than 14 mag. To avoid additional work of transforming different photometry systems to an unified one for all LAMOST spectra, we use 2MASS (Cutri et al. 2003), which is in near infrared bands covering most of the magnitude range but with less accuracy, as a compromising solution to provide a homogeneous distance estimation for both bright and faint stars.

We use a dense grid of synthetic isochrones (Girardi et al. 2002; Marigo et al. 2008) at a typical age of 10 Gyr, which is suitable for the halo stars. The library is tailored to contain only RGB+SGB tracks. The posterior probability density function (PDF) of the absolute magnitude in  $K$  band,  $M_K$ , for a giant star can be obtained from the following equation according to Bayes' theorem:

$$p(M_K|(J-K)_0, [\text{Fe}/\text{H}], \text{isochrone}) = p((J-K)_0, [\text{Fe}/\text{H}]|M_K, \text{isochrone})p(M_K), \quad (2)$$

where the term  $p((J-K)_0, [\text{Fe}/\text{H}]|M_K, \text{isochrone})$  is the likelihood and the prior  $p(M_K)$  the luminosity function.  $(J-K)_0$  is the reddening corrected color index. We adopt the extinction map from Schlegel, Finkbeiner, & Davis (1998) and the extinction parameters for  $J$  and  $K$  bands from Cardelli, Clayton, & Mathis (1989) given  $R_V = 3.1$ . In this work, we adopt the luminosity function at 10 Gyr from the library website<sup>5</sup>, which is based on Chabrier IMF (Chabrier 2001). Note that  $(J-K)_0$  is estimated from photometric observation and  $[\text{Fe}/\text{H}]$  is derived from spectroscopic survey data, therefore, they are supposed to be two independent measurements. Then the likelihood can be separated into:

$$p((J-K)_0, [\text{Fe}/\text{H}]|M_K, \text{isochrone}) = p((J-K)_0|M_K, \text{isochrone})p([\text{Fe}/\text{H}]|M_K, \text{isochrone}). \quad (3)$$

The two terms on the right hand side of equation (3) can be specified as the following equations providing Gaussian errors for the measurement of  $J-K$  and  $[\text{Fe}/\text{H}]$ .

$$p((J-K)_0|M_K, \text{isochrone}) \sim \exp\left(-\frac{((J-K)_{iso} - (J-K)_0)^2}{2\sigma_{J-K}^2}\right), \quad (4)$$

$$p([\text{Fe}/\text{H}]|M_K, \text{isochrone}) \sim \exp\left(-\frac{([\text{Fe}/\text{H}]_{iso} - [\text{Fe}/\text{H}])^2}{2\sigma_{[\text{Fe}/\text{H}]}^2}\right), \quad (5)$$

where the variables in equations (4) and (5) with subscript *iso* are quantities from isochrone and those without subscript are from the observed data of a given star. The  $\sigma_{J-K}$  and  $\sigma_{[\text{Fe}/\text{H}]}$  are measurement uncertainties of  $(J-K)_0$  and  $[\text{Fe}/\text{H}]$  for the star.

The posterior PDF of the absolute magnitude for a star derived from equation (2) is then converted to that of DM by applying apparent  $K$  band magnitude. In the next sections we use the median value of DM as the best estimated value. The  $-\sigma$  and  $+\sigma$  of DM are defined by the 15% and 85% percentiles of the PDF.

## 5.2. Calibration with globular clusters

Distance estimation using synthetic isochrones, as described above, needs to be calibrated observationally. A usual practice is to use globular clusters (GCs). The main issue is that the isochrones in the observed data are different from that of the synthetic spectra, given the same metallicity, and the offset becomes larger for metal-poor objects (see figure 11). We select 7 GCs (see table 3), which are bright and cover a wide range in metallicity, to calibrate the color offsets. The giant members of the GCs are manually selected from the  $(J-K)_0$  vs.  $K$  diagram of stars within a small radius (defined in the 5th column of table 3) around the center of each cluster (the red dots shown in the color-magnitude diagrams in figures 11 and 12). Note that although they are not definite GC members, the contaminations do not affect the result since we replace the individual metallicities for the selected stars with the mean value of the GC in the calibration and later test.

We apply a polynomial surface model to obtain the corrected offset color index,  $\Delta(J-K)$ , as a function of  $(J-K)_0$  and  $[\text{Fe}/\text{H}]$ . The best fit is

<sup>5</sup><http://stev.oapd.inaf.it/cgi-bin/cmd>



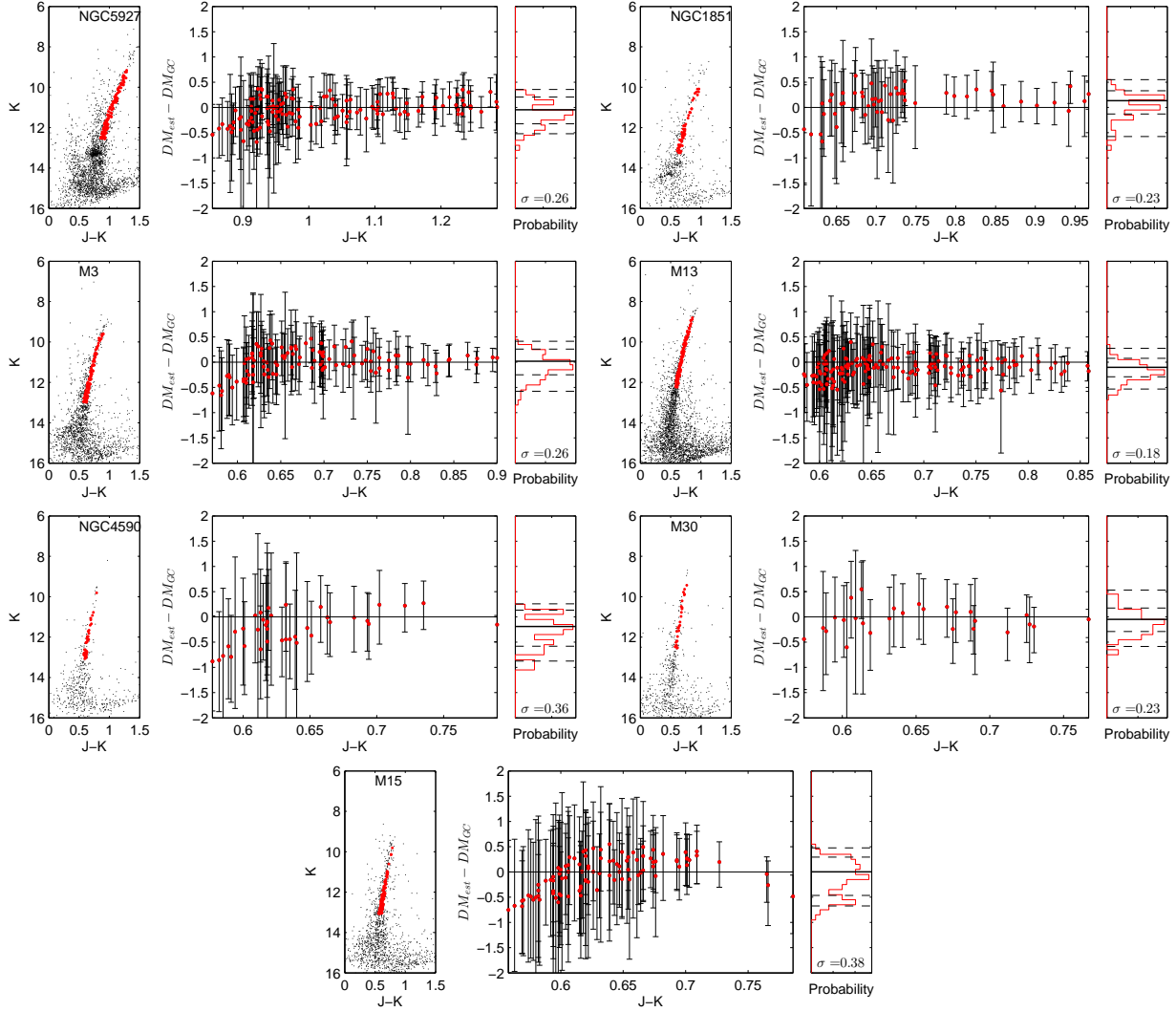


Fig. 12.— The distance performance tested with globular clusters NGC5927, NGC1851, M3, M13, NGC4590, M30, and M15. For each panel, the left plot shows the color-magnitude diagram for the stars around the cluster (black dots) and the selected K giant members (red dots). The middle one shows the residuals of the distance moduli of the K giant members (red dots) with  $1\text{-}\sigma$  uncertainties as a function of  $J - K$ . The right one shows the distribution of the residual DM (red lines). The solid and dashed horizontal lines indicate the peak,  $1\text{-}\sigma$ , and  $2\text{-}\sigma$  values.

Table 3: The parameters of 7 GCs for calibration

GC names	Distance <sup>a</sup>	[Fe/H] <sup>a</sup>	E(B-V) <sup>b</sup>	selection radius <sup>c</sup>
	(kpc)	dex	mag	degree
NGC5927	7.7	-0.49	0.45	0.1
NGC1851	12.1	-1.18	0.02	0.07
M3	10.2	-1.5	0.02	0.2
M13	7.1	-1.53	0.02	0.32
NGC4590	10.3	-2.23	0.05	0.07
M30	8.1	-2.27	0.03	0.07
M15	10.4	-2.37	0.1	0.32

<sup>a</sup>Harris (1996, revision 2010)

<sup>b</sup>Color excess in  $B - V$ , adopting the values from Harris (1996, revision 2010)

<sup>c</sup>The radius used to select the giant members around the center of each GC.

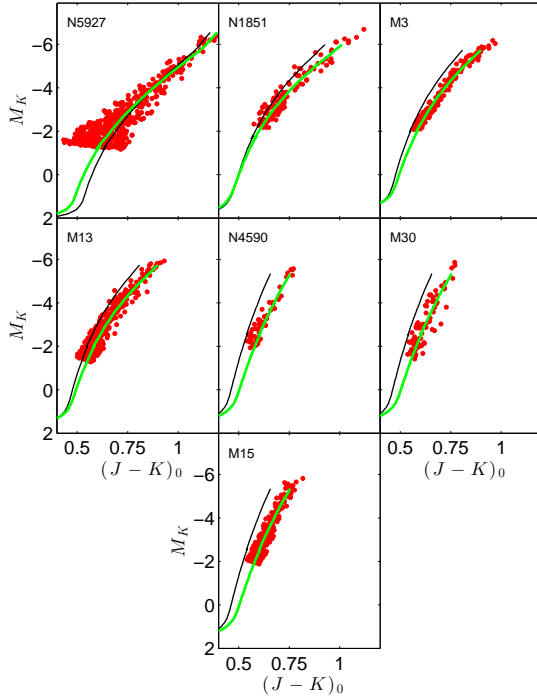


Fig. 11.— The black lines are the synthetic isochrones with the same  $[\text{Fe}/\text{H}]$  as the globular clusters. The thick green lines are the calibrated isochrones according to equation (6). The red dots are the K giant members of the globular clusters.

$$\begin{aligned} \Delta(J-K) = & -0.1421 \pm 0.0187 + \\ & 0.08454 \pm 0.0209(J-K)_0 - \\ & 0.06482 \pm 0.0170[\text{Fe}/\text{H}] - \\ & 0.09195 \pm 0.0166(J-K)_0[\text{Fe}/\text{H}] - \\ & 0.0194 \pm 0.0026[\text{Fe}/\text{H}]^2. \end{aligned} \quad (6)$$

Figure 11 demonstrates the performance of the calibration. The red dots are the K giant members of the 7 globular clusters, the black lines show the synthetic isochrones with same metallicity as the corresponding GCs. After employing the calibration of  $(J-K)_0$  using equation (6), the isochrones are shifted to the green lines, well fit the GC data.

Adding  $\Delta(J-K)$  into equation (4), we obtain the new likelihood of  $(J-K)_0$ :

$$\begin{aligned} p((J-K)_0 | M_K, \text{isochrone}) \sim \\ \exp\left(-\frac{((J-K)_{iso} + \Delta(J-K) - (J-K)_0)^2}{2\sigma_{J-K}^2}\right). \end{aligned} \quad (7)$$

The distances of the K giant members are derived by using the above calibration. Figure 12 shows the performance of the estimation using the member giant stars of the 7 GCs. Figure 13 shows the errors of DM (denoted as  $\sigma_{DM}$ ) as a function of apparent magnitude  $K$ . At 11 mag, the median value of the full K giant sample, the mean error of DM is about 0.6 mag, or  $\sim 30\%$  in distance. This is a factor of 2.5 larger than that of Xue et al. (2014). The use of 2MASS, which is in infrared, shallower, and less accurate than SDSS, is likely responsible for the poorer performance. The distance measurement is inevitably compromised until a unified, high accuracy optical photometric catalog is established for the LAMOST survey.

### 5.3. Measuring distance for LAMOST K giant stars

Before deriving distance of the K giant stars, the sample is purified by the following procedure: i) only the later/latest epoch observation is kept for duplicated objects; ii) the non-stellar targets are removed; iii) only the stars with  $5 < K < 15.5$  mag and the  $K$  magnitude error lower than 0.1 mag are selected; iv) the stars with  $|b| < 10^\circ$

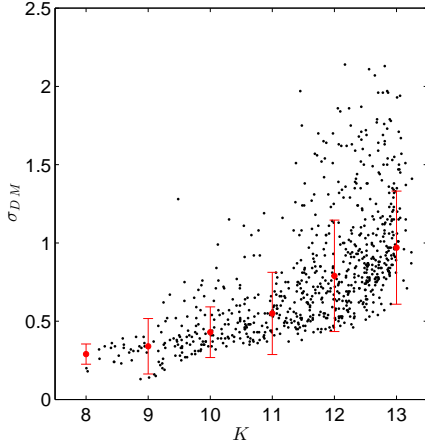


Fig. 13.— The black dots are the error of DM,  $\sigma_{DM}$ , of the giant members of the 7 globular clusters as a function of apparent magnitude  $K$ . The red dots present the median values at given  $K$  magnitude.

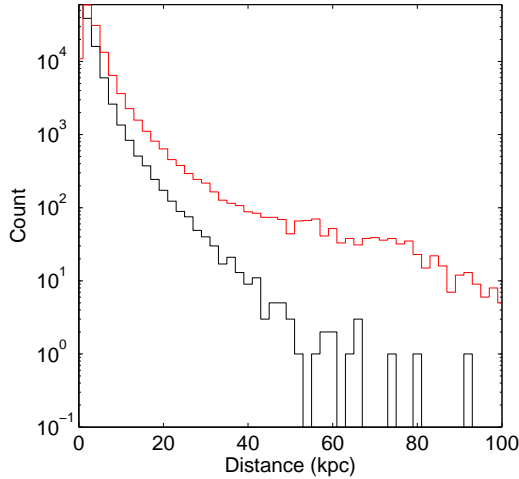


Fig. 14.— The distribution of the distance for the K giant stars. The red and black lines represent for the K giant stars selected from the SVM classifier and from the LAMOST derived parameters, respectively.

are removed to avoid the high extinction and the majority of the disk K giant stars; v) the K giant stars with  $[\text{Fe}/\text{H}]_{LM2D} < -3$  dex, which may not be reliable, are also excluded; and vi) the red clump stars selected according to Appendix A are excluded, since the distance estimation method described in section 5 is not suitable for them. After these cuts, there are 134,597 RGB/SGB-like K giant remaining.

Meanwhile, we clip the *true* K giant samples selected from LAMOST parameter catalog (see subsection 3.5) using the similar criteria mentioned in previous paragraph, and obtain 76,822 *true* K giant stars left for distance estimation with  $[\text{Fe}/\text{H}]_{LM}$ .

In figure 14, the spatial distribution of the K giant stars based on  $[\text{Fe}/\text{H}]_{LM2D}$  (red) and LAMOST  $[\text{Fe}/\text{H}]_{LM}$  (black) are plotted. The star count of LM2D at 20 kpc is a factor of 4 larger than that of LAMOST sample at the same distance. The factor increases to around 10 when the distance is larger than 40 kpc. Therefore, our identification method significantly increases the K giant sampling of the survey, which is more efficient in larger distances, allowing the search for the kinematic substructure in a much larger volume.

#### 5.4. Caveats of the distance estimation

Because the accuracy of distance estimation depends on the accuracy of the metallicity, the systematic bias of the estimated metallicity, as shown in figure 9, may induce a systematic bias in estimated distance. Figure 15 shows the distribution of the residual DM for the K giant members of M3. With the true metallicity, i.e., -1.5 dex according to table 3, the estimated distance does not show any significant systematic bias. However, when the metallicity of the member stars is overestimated by 0.15 (0.3) dex, i.e.,  $[\text{Fe}/\text{H}] = -1.35$  (-1.2 dex), the median residual DM shifts 0.15 (0.35) mag towards left, which equivalent with 7% (17%) underestimation of the distance.

For the halo and thick disk K giant stars, the distance estimation based on isochrones with age of 10 Gyr should be acceptable. However, for younger stars in the disk this may lead to a systematic bias in the distance estimation. Figure 16 shows that for the metal-rich globular cluster NGC5927, which true metallicity is -0.49, the

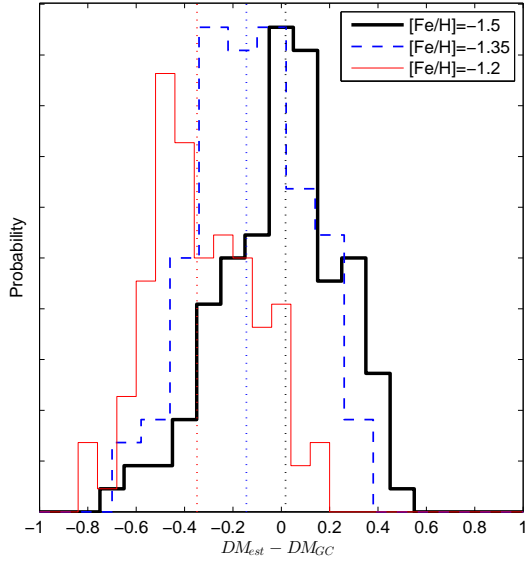


Fig. 15.— The effect of the overestimation of the metallicity in the estimated distance. The histograms show the residual distribution of the distance moduli of the M3 member K giant stars when the applied metallicity is -1.5 (the true metallicity of the globular cluster, shown in thick black line), -1.35 (blue dashed line), and -1.2 (red thin line), respectively. Their median residuals are marked using vertical black, blue, and red dot lines, respectively.

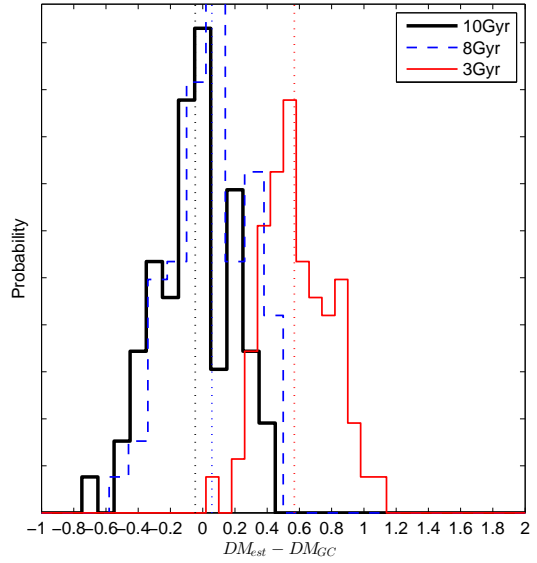


Fig. 16.— The effect of the age in the estimated distance. The histograms show the residual distribution of the distance moduli of the NGC5927 member K giant stars when the applied age of the isochrones is 10 (thick black line), 8 (blue dashed line), and 3 Gyr (red thin line), respectively. Their median residuals are marked using vertical black, blue, and red dot lines, respectively.

residual of the DM shifts by about 0.6 mag toward right when a set of 3 Gyr isochrones are applied in the distance estimation. In other words, the distance estimation method may underestimate by  $\sim 30\%$  for this case.

As a summary, the distance may be underestimated either for metal-poor stars since both the LAMOST pipeline and the LM2D method tend to slightly overestimate their metallicity or for young disk stars due to the application of relatively old isochrones.

## 6. The Sagittarius stream

In order to assess and demonstrate the validity of our catalog, we carry out the following scientific verification. The Milky Way's stellar halo is known to host many tidal streams and substructures, which can be identified using catalogs such as the one presented here. The most prominent stream originates from the Sagittarius dwarf galaxy (Ibata et al. 2001). This stream has been detected by a number of authors using a variety of datasets, both photometric and spectroscopic (Newberg et al. 2002; Majewski et al. 2003; Belokurov et al. 2006; Dohm-Palmer et al. 2001; Majewski et al. 2004a; Monaco et al. 2007; Carlin et al. 2012, etc.). These detections can then be used to probe the potential of the Milky Way (e.g., Helmi & White 2001; Law, Johnston, & Majewski 2005; Fellheuer et al. 2006; Law & Majewski 2010) or to estimate the properties of the host galaxy, which was originally much more massive than it is today (Niederste-Ostholt et al. 2010).

To aid our search for Sagittarius members, we convert the equatorial coordinates into the Belokurov et al. (2014) system (itself based on the system of Majewski et al. 2003), described by two angles  $\tilde{\Lambda}_\odot$  and  $\tilde{B}_\odot$ . This system has the stream located along the equator, with the core at  $\tilde{\Lambda}_\odot = 0$  and this angle increasing in the direction of motion, i.e. the start of the leading stream lies at positive  $\tilde{\Lambda}_\odot$  and the trailing stream at negative  $\tilde{\Lambda}_\odot$ . An illustration of this coordinate system is shown in figure 2 of Belokurov et al. (2014), where it can be seen that the stream lies in the region  $|\tilde{B}_\odot| \lesssim 10$  deg. We correct for extinction using the maps of Schlegel, Finkbeiner, & Davis (1998).

We begin by searching for stream members in the South Galactic cap region, where one can easily detect material from the trailing arm. This has been studied in detail and its properties well known (see Koposov et al. 2012, for an overview). The distances at this location are around 30 kpc and the velocities significantly offset from the background halo at around  $-150$  to  $-100$  km/s. This can be seen clearly in the upper panel of figure 17, with the stream members separated from the bulk of the halo population. In the middle panel we present a histogram of radial velocities for all stars with  $|\tilde{B}_\odot| < 15$  deg and, for comparison, show the distribution of particles in this region from the model of Law & Majewski (2010), where we have only included stars which have been stripped in the past 3 Gyr. There is good agreement here, which is unsurprising as the model was tuned to match the velocity signal in this part of the stream, but it is reassuring that our K giant sample is able to pick up a clean sample of Sagittarius members. The lower panel shows the distances for these stars, retaining only those which have  $-160 < v_{gst} < -100$  km/s. Here we can see that there is a slight offset in distance between our data and Law's model, at a level of around 20 per cent. Again it is known that the Law model does not have problems in this region, as evidenced by good agreement with main-sequence turn-off stars from Koposov et al. (2012). Therefore this offset is most likely due to systematics in our distances, which we discuss below.

Another avenue to verify the catalogue is to look for material at the trailing tail's apo-centre. This was recently detected by Belokurov et al. (2014) and Drake et al. (2013) in both blue horizontal branch stars and red giants and lies at a distance of around 100 kpc. Although the radial velocity is, by definition, close to zero in this region, the signal should be relatively strong due to the small number of smooth halo stars this far out in the halo.

This indeed proves to be the case, as can be seen in figure 18. The top panel shows how the radial velocity varies as a function of Galacto-centric distance for all stars in our catalogue and there is already a hint of a detection in the clump of stars around 80 kpc. The nature of these stars is uncovered in the middle panel, where we show their location on the sky. They are clearly not

drawn uniformly from the underlying population, but are obviously clumped around  $\tilde{B}_\odot \approx 0$  deg and  $\tilde{\Lambda}_\odot \approx 170$  deg, which is precisely where the Belokurov detection lies. This finding is reinforced in the lower-panel, which shows the distribution of these stars as a function of  $\tilde{\Lambda}_\odot$  - their distribution is dramatically different from the general footprint of the LAMOST survey and so this finding cannot be explained by a quirk of the survey strategy.

We compare our detection to that of Belokurov et al. (2014) in figure 19. The top panel shows the distribution of helio-centric distances for our K giants in the region of Belokurov’s apo-centre detection ( $|\tilde{B}_\odot| < 30$  deg,  $160 < \tilde{\Lambda}_\odot < 180$  deg). This figure demonstrates that the bulk of these stars are in good agreement with their velocities (middle panel), although our distances again appear to be systematically offset when compared to Belokurov’s blue horizontal branch star distances (bottom panel). As before, this distance offset is at around the 20 per cent level.

To summarize, we have shown that our catalogue is well-suited to the detection of halo substructures and we believe that it will be a rich resource for finding further streams. Such endeavors are beyond the scope of this work, but we are now pursuing this goal and will report our findings in a future publication. However, our distance estimates do need to be handled with care. Systematic offsets could be due to various factors, such as incorrectly estimated metallicities or ages (Section 5.4), or possibly contamination from red clump stars (Appendix A).

## 7. Conclusions

We have established a SVM classifier directly from the spectra features, and then apply it to LAMOST spectra for K giant star selection. The method does not depend on the stellar parameters, e.g.,  $T_{\text{eff}}$  and  $\log g$ , thus has a broader range of capability to work on spectra with S/N as low as 3. Tested with SDSS, MILES, and LAMOST data, the SVM classifier can select K giant stars from the survey dataset with 70-80% completeness and a few percent contamination. From the DR1 released  $\sim 1.9$  million stellar spectra, we identified about 290,000 K giant stars. Consequently, we expect that, when the survey will be concluded in five years, there will be a factor of 4 more K giant stars to be observed.

In order to estimate the distance of the K giant stars, we have firstly estimated the metallicity using LM2D method. Comparisons with SDSS parameters indicates that the total error of the estimation is between 0.1 dex and 0.3 dex. The advantage of the estimation method is that we can provide metallicity for the identified K giant stars with signal-to-noise ratio down to 3.

We have then developed a Bayesian method to estimate the distance of the K giant stars using 2MASS photometry and the estimated metallicity from LM2D. The synthetic isochrone-based method is calibrated with 7 globular clusters. Therefore, the systematic bias due to the discrepancy between the synthetic and observed data is corrected. The uncertainty of the distance estimation is investigated using the same globular clusters, which covers a wide range in metallicities. We conclude that the uncertainty in DM is around 0.5 mag at  $K \sim 11$  mag, corresponding to about 30% in distance.

Given the distance and radial velocities of the K giant stars selected from the survey, we successfully identified many candidate members of the Sgr stream. These identifications demonstrate that there may be thousands of K giant members of such tidal substructures to be observed and identified over a broad area of the sky by the end of the five-year survey. Consequently, it will be one of the most important homogeneous spectroscopic dataset to map the kinematics as well as the chemical abundance of them. This will provide a tight constraint on the dark matter mass in the Galactic

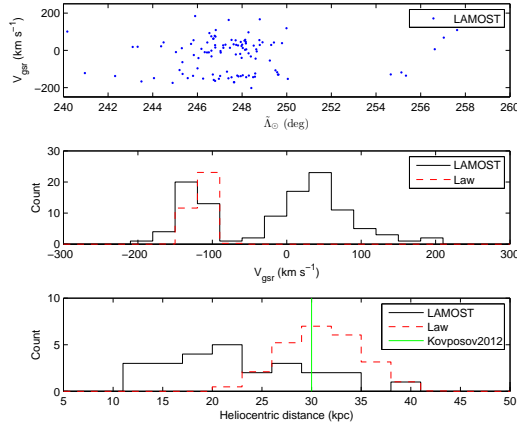


Fig. 17.— The detection of Sagittarius trailing debris. The top panel shows the velocity signature of the stream in our K giant sample around  $|\tilde{B}_\odot| < 15$  deg. The detection is even more evident when we plot a histogram of Galacto-centric radial velocities ( $V_{\text{gsr}}$ ) for these stars (middle panel), as can be seen from the peak around  $-160$  to  $-100$  km/s. The dashed line shows the prediction from the model of Law & Majewski (2010), where we have only included stars stripped in the past 3 Gyr. The bottom panel shows the distances of these giant stars with velocities consistent with the stream (i.e.  $-160 < V_{\text{gsr}} < -100$  km/s), along with the prediction from Law’s model. There is a small ( $\sim 20$  per cent) offset, which is discussed in Section 6.

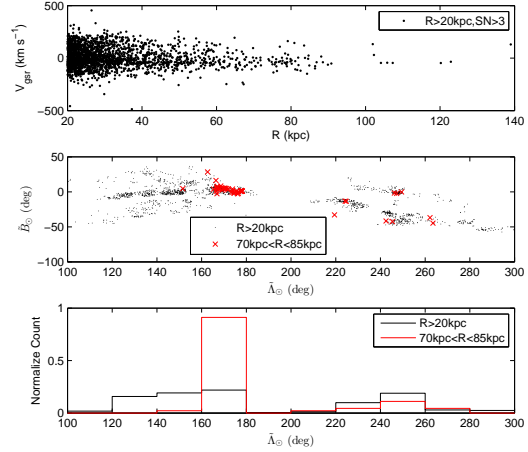


Fig. 18.— The detection of distant Sagittarius debris, corresponding to the apo-centre of the trailing tail. The top panel shows the distribution of galacto-centric radial velocity ( $V_{\text{gsr}}$ ) and galacto-centric distance ( $R$ ) for all stars with signal-to-noise larger than 3. There appears to be a clump of stars around 80kpc. The spatial distribution of these distant stars ( $70 < R < 85$  kpc; red crosses) is compared with the underlying distribution ( $R > 20$  kpc; black points). The middle panel shows that these stars are also concentrated in a small range of  $\tilde{A}_\odot$  between  $160^\circ$  and  $180^\circ$ . The bottom panel shows the histograms of  $\tilde{A}_\odot$  for these two distributions and, again, it is clear that the distant giants are not drawn uniformly from the underlying distribution.

halo as well as the forming history of the substructure themselves.

This work is supported by the Strategic Priority Research Program "The Emergence of Cosmological Structures" of the Chinese Academy of Sciences, Grant No. XDB09000000 and the National Key Basic Research Program of China 2014CB845700. CL acknowledge the National Science Foundation of China under grants 11373032 and U1231119. JLC and HJN acknowledge National Science Foundation under grant AST 09-37523. XXX acknowledges the Alexandra Von Humboldt foundation for a fellowship and the National Natural Science Foundation of China under grants 11103031, 11233004 and 11003017.

Guoshoujing Telescope (the Large Sky Area Multi-Object Fiber Spectroscopic Telescope LAMOST) is a National Major Scientific Project built by the Chinese Academy of Sciences. Funding for the project has been provided by the National Development and Reform Commission. LAMOST is operated and managed by the National Astronomical Observatories, Chinese Academy of Sciences.

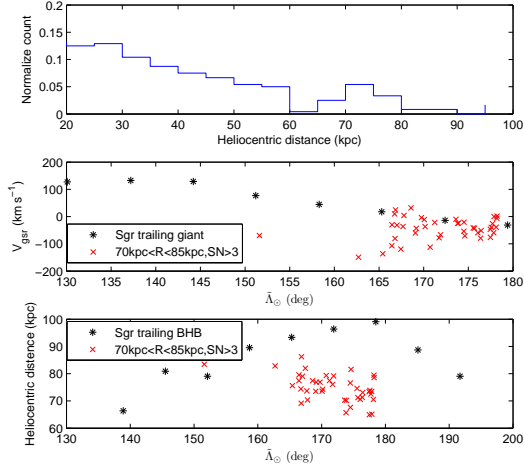


Fig. 19.— A comparison between our detection of the distant Sagittarius debris and that of Belokurov et al. (2014). The top panel shows the distribution of helio-centric distances for our LAMOST giants (with  $S/N > 3$ ) in the region where these distant stars lie ( $|\bar{B}_\odot| < 30$  deg,  $160 < \bar{A}_\odot < 180$  deg). A clear overdensity is visible at helio-centric distances of around 65 to 80 kpc, corresponding to Galacto-centric distances  $70 < R < 85$  kpc. The middle panel compares our LAMOST giant velocities (red crosses) to Belokurov’s SDSS giant velocities (black asterisks). The same notation is used in the lower panel, where we compare our giant distances to their blue horizontal branch star distances. There is a small ( $\sim 20$  per cent) offset, which is discussed in Section 6.



### A. Remove the red clump stars

Even in low extinction regions, there are quite a lot red clump stars in the K giant sample. The distance estimation method is not suitable for this type of stars since the isochrone is only for RGB/SGB stars. Hence, we have to identify and remove them from the samples so that we can use a purified dataset in the study of the spatial distribution of the K giant stars.

The top-left panel of figure 20 shows the location of the red clump stars in  $T_{\text{eff}}$  vs.  $\log g$  plane. We use a polygon (red lines) to select the possible red clump stars from their  $T_{\text{eff}}$  and  $\log g$ . However, for the majority of the identified K giant stars, there is no measurement of  $T_{\text{eff}}$  and  $\log g$ . We turn to use alternative quantities directly measured from the spectra. The top-right panel shows the  $EW_{Mgb}$  vs.  $[\text{Fe}/\text{H}]_{LM2D}$  for the selected K giant stars (black filled contours). The possible red clump stars inside the red polygon are mostly concentrated into a narrower region (yellow contours). For simplicity, we use the thick blue polygon shown in the top-right panel to locate the red clump stars. Subsequently, we remove all the stars that fall in the polygon to exclude the red clump star contamination. After the exclusion of the possible red clump stars, the  $T_{\text{eff}}$  vs.  $\log g$  distribution of the rest of the K giant stars is shown in the bottom-left panel. The clump centered at  $T_{\text{eff}} \sim 4900$  and  $\log g \sim 2.5$  in the top-left panel is now almost disappeared. On the other hand, the excluded stars located within the thick blue polygon are concentrated as a clump in the  $T_{\text{eff}}$  vs.  $\log g$  plane, as shown in the bottom-right panel.

A quantitative assessment of the performance of the red clump star exclusion method is very difficult since we cannot definitely determine whether an individual field star is a red clump star. However, figure 20 shows that, qualitatively, the simple polygon exclusion works quite well in the removing of the red clump stars. This is sufficient for the spatial overview of the K giant stars, losing a small fraction of K giant stars being mistakenly classified as red clump stars and keeping a smaller fraction of contaminated red clump stars in the sample. The distances for the contaminated red clump stars are likely underestimated because that the RGB/SGB stars with same  $J - K$  and  $[\text{Fe}/\text{H}]$  have fainter absolute magnitude than that of the red clump stars. An elegant classification of the red clump stars deserves an another specific project and beyond the scope of this work.

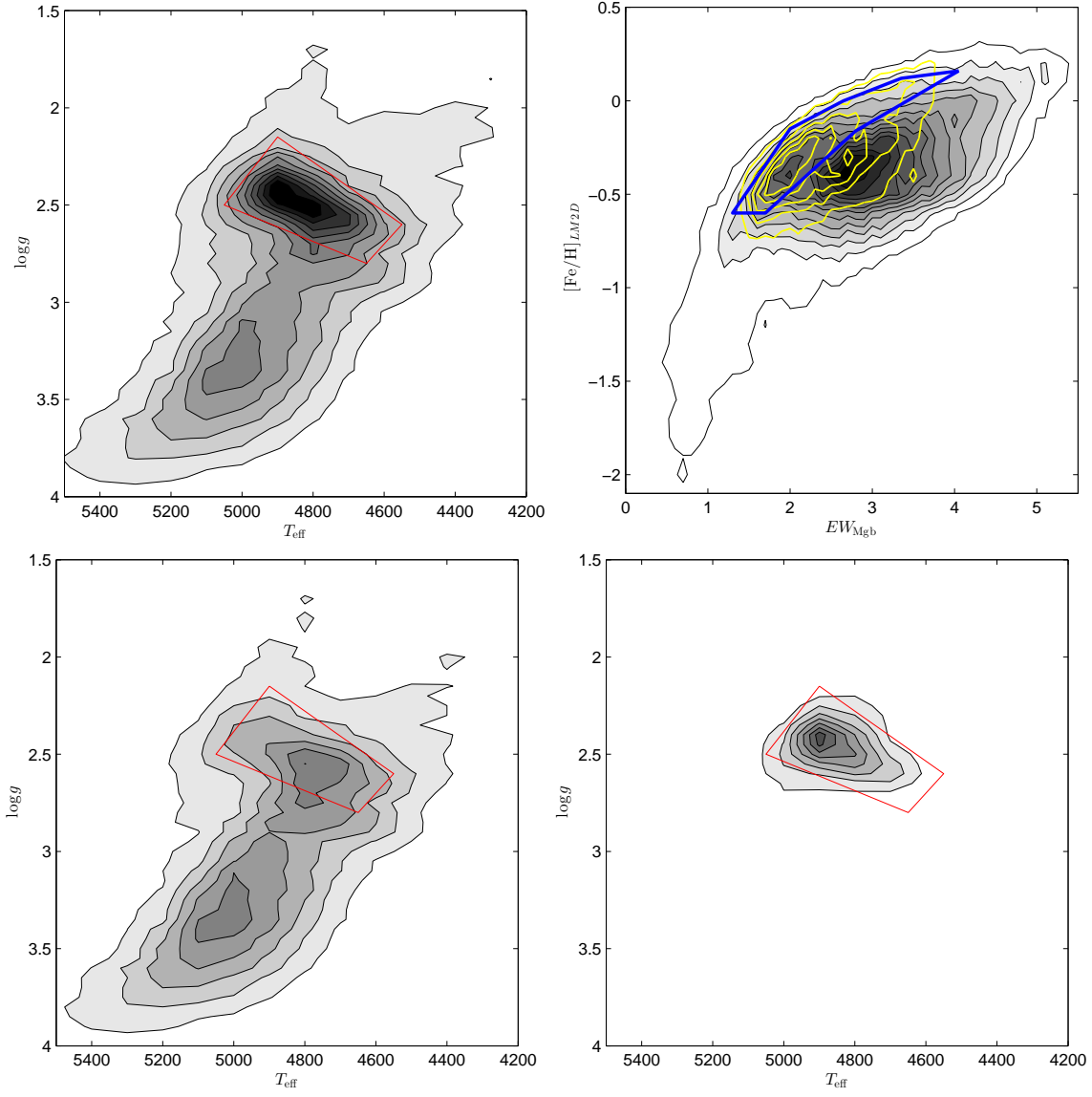


Fig. 20.— *Top-left panel:* The contour shows the distribution of the stars with LAMOST stellar parameters in  $T_{\text{eff}}$  vs.  $\log g$ . The red polygon indicates the location of the red clump stars. *Top-right panel:* The K giant stars (black filled contours) in  $EW_{\text{Mgb}}$  vs.  $[\text{Fe}/\text{H}]_{\text{LM2D}}$ . The overlapped yellow contours are the stars located in the red polygon. The levels of the both contours are the star counts per unit  $EW_{\text{Mgb}} \times [\text{Fe}/\text{H}]_{\text{LM2D}}$  varying from 5000 to 80000 with step size of 5000. The thick blue polygon is the simple selection criteria to define the red clump stars in  $EW_{\text{Mgb}}$  vs.  $[\text{Fe}/\text{H}]_{\text{LM2D}}$ . *Bottom-left panel:* The  $T_{\text{eff}}$  vs.  $\log g$  distribution of K giant stars excluded all samples inside the thick blue polygon defined in the top-right panel. *Bottom-right panel:* The  $T_{\text{eff}}$  vs.  $\log g$  distribution of the stars inside the thick blue polygon, in which most of the red clump stars are located.

## REFERENCES

- Ahn, C. P., Alexandroff, R., Allende Prieto, C., Anderson, S. F., Anderton, T., Andrews, B. H., Aubourg, É., et al., 2012, *ApJS*, 203, 21
- Bailer-Jones, C. A. L., Smith, K. W., Tiede, C., Sordo, R., Vallenari, A., 2008, *MNRAS*, 391, 1838
- Bailer-Jones, C. A. L., Andrae, R., Arcay, B., Astraatmadja, T., Bellas-Velidis, I., Berihuete, A., Bijaoui, A., Carrión, C. et al., 2013, arXiv:1309.2157, accepted to *A&A*
- Belokurov, V., Zucker, D. B., Evans, N. W., Gilmore, G., Vidrih, S., Bramich, D. M., Newberg, H. J., et al., 2006, *ApJ*, 642, 137
- Belokurov, V., Evans, N. W., Bell, E. F., Irwin, M. J., Hewett, P. C., Koposov, S., Rockosi, C. M., et al., 2007, *ApJ*, 657, 89
- Belokurov, V., Koposov, S. E., Evans, N. W., Peñarrubia, J., Irwin, M. J., Smith, M. C., Lewis, G. F., Gieles, M., et al. 2014, *MNRAS*, 437, 116
- Bonaca, A., Geha, M., Kallivayalil, N., 2012, *ApJ*, 760, 6
- Cardelli, J. A., Clayton, G. C., & Mathis, J. S., 1989, *ApJ*, 345, 245
- Carlin, J. L., Majewski, S. R., Casetti-Dinescu, D. I., Law, D. R., Girard, T. M., Patterson, R. J., 2012, *ApJ*, 744, 25
- Chabrier, G., 2001, *ApJ*, 554, 1274
- Cortes C., Vapnik V., 1995, *Machine Learning*, 20, 273
- Cui, X., Zhao, Y., Chu, Y., Li, G., Li, Q., Zhang, L., Su, H., et al., 2012, *RAA*, 12, 1197
- Cutri, R. M., Skrutskie, M. F., van Dyk, S., Beichman, C. A., Carpenter, J. M., Chester, T., Cambresy, L., 2003, *yCat*, 2246, 0
- Deng, L., Newberg, H. J., Liu, C., Carlin, J. L., Beers, T. C., Chen, L., Chen, Y., et al., 2012, *RAA*, 12, 777
- Dohm-Palmer, R. C., Helmi, A., Morrison, H., Mateo, M., Olszewski, E. W., Harding, P., Freeman, K. C., et al., 2001, *ApJ*, 555, 37
- Drake, A. J., Catelan, M., Djorgovski, S. G., Torrealba, G., Graham, M. J., Mahabal, A., Prieto, J. L., Donalek, C., et al., 2013, *ApJ*, 765, 154
- Fellhauer, M., Belokurov, V., Evans, N. W., Wilkinson, M. I., Zucker, D. B., Gilmore, G., Irwin, M. J., et al., 2006, *ApJ*, 651, 167
- Girardi, L., Bertelli, G., Bressan, A., Chiosi, C., Groenewegen, M. A. T., Marigo, P., Salasnich, B., Weiss, A., 2002, *A & A*, 391, 195
- Grillmair, C. J., Dionatos, O., 2006, *ApJ*, 643, L17
- Harris, W. E. 1996, *AJ*, 112, 1487
- Helmi, A., & White, S. D. M. 2001, *MNRAS*, 323, 529
- Ibata, R., Irwin, M., Lewis, G. F., Stolte, A., 2001, *ApJ*, 547, 133
- Klypin, A., Kravtsov, A. V., Valenzuela, O., Prada, F. 1999, *ApJ*, 522, 82
- Koposov, S., Belokurov, V., Evans, N. W., Hewett, P. C., Irwin, M. J., Gilmore, G., Zucker, D. B., et al., 2008, *ApJ*, 686, 279
- Koposov, S. E., Rix, H.-W., Hogg, D. W., 2010, *ApJ*, 712, 260
- Koposov, S. E., Belokurov, V., Evans, N. W., Gilmore, G., Gieles, M., Irwin, M. J., Lewis, G. F., Niederste-Ostholt, M. et al., 2012, *ApJ*, 750, 1
- Law, D. R., Johnston, K. V., Majewski, S. R., 2005, *ApJ*, 619, 807
- Law, D. R., Majewski, S. R., 2010, *ApJ*, 714, 229
- Liu, C., Bailer-Jones, C. A. L., Sordo, R., Vallenari, A., Borrachero, R., Luri, X., Sartoretti, P., 2012, *MNRAS*, 426, 2463
- Liu, X.-W., Yuan, H.-B., Huo, Z.-Y., Deng, L.-C., Hou, J.-L., Zhao, Y.-H., Zhao, G., Shi, J.-R., et al., 2013, arXiv:1306.5376, to appear in *Proceedings of the IAUS 298*
- Luo, A., Zhang, H., Zhao, Y., Zhao, G., Cui, X., Li, G., Chu, Y., et al., 2012, *RAA*, 12, 1243
- Majewski, S. R., Skrutskie, M. F., Weinberg, M. D., Ostheimer, J. C. 2003, *ApJ*, 599, 1082

- Majewski, S. R., Kunkel, W. E., Law, D. R., Patterson, R. J., Polak, A. A., Rocha-Pinto, H. J., Crane, J. D., et al., 2004a, *AJ*, 128, 245
- Majewski, S. R., Ostheimer, J. C., Rocha-Pinto, H. J., Patterson, R. J., Guhathakurta, P., Reitzel, D., 2004b, *ApJ*, 615, 738
- Marigo, P., Girardi, L., Bressan, A., Groenewegen, M. A. T., Silva, L., Granato, G. L., 2008, *A & A*, 482, 883
- Martin, C., Carlin, J. L., Newberg, H. J., Grillmair, C., 2013, *ApJ*, 765, L39
- Monaco, L.; Bellazzini, M.; Bonifacio, P.; Buzzoni, A.; Ferraro, F. R.; Marconi, G.; Sbordone, L.; Zaggia, S., 2007, *A & A*, 464, 201
- Newberg, H. J., Yanny, B., Rockosi, C., Grebel, E. K., Rix, H.-W., Brinkmann, J., Csabai, I., et al., 2002, *ApJ*, 659, 245
- Newberg, H. J., Yanny, B., Cole, N., Beers, T. C., Re Fiorentin, P., Schneider, D. P., Wilhelm, R., 2007, *ApJ*, 668, 221
- Newberg, H. J., Yanny, B., Willett, B. A., 2009, *ApJ*, 700, 61
- Niederste-Ostholt, M., Belokurov, V., Evans, N. W., Peñarrubia, J., 2010, *ApJ*, 712, 516
- Pietrinferni, A., Cassisi, S., Salaris, M., & Castelli, F., 2004, *ApJ*, 612, 168
- Prugniel, P., Koleva, M., Ocvirk, P., Le Borgne, D., Soubiran, C. 2007, in *IAU Symposium*, vol. 241, eds. A. Vazdekis, R. F. Peletier, 68
- Rocha-Pinto H. J., Majewski S. R., Skrutskie M. F., Crane J. D., Patterson R. J. 2004, *ApJ*, 615, 732
- Saglia, R. P., Tonry, J. L., Bender, R., Greisel, N., Seitz, S., Senger, R., Snigula, J., et al., 2012, *ApJ*, 746, 128
- Sánchez-Blázquez, P., Peletier, R. F., Jiménez-Vicente, J., Cardiel, N., Cenarro, A. J., Falcón-Barroso, J., Gorgas, J., et al., 2006, *MNRAS*, 371, 703
- Schlegel, David J.; Finkbeiner, Douglas P.; Davis, Marc, 1998, *ApJ*, 500, 525
- Shi, W. B., Chen, Y. Q., Carrell, K., Zhao, G., 2012, *ApJ*, 751, 130
- Smith, K. W., Bailer-Jones, C. A. L., Klement, R. J., Xue, X. X., 2010, *A & A*, 522, 88
- Martin, C., Carlin, J. L., Newberg, H. J., Grillmair, C., 2013, *ApJ*, 765, 39
- Vivas, A. K., Zinn, R., 2006, *AJ*, 132, 714
- Willett, B. A., Newberg, H. J., Zhang, H., Yanny, B., Beers, T. C., 2009, *ApJ*, 697, 207
- Worthey, G., Faber, S. M., Gonzalez, J. J., Burstein, D., 1994, *ApJS*, 94, 687
- Wu, Y., Luo, A., Li, H., Shi, J., Prugniel, P., Liang, Y., Zhao, Y., et al., 2011, *RAA*, 11, 924
- Xue, X.-X., Ma, Z., Rix, H. W., Morrison, H. L., Harding, P., Beers, T. C., Ivans, I. I., et al., 2014, *ApJ*, 784, 170
- Yanny, B., Newberg, H. J., Grebel, E. K., Kent, S., Odenkirchen, M., Rockosi, C. M., Schlegel, D., et al., 2003, *ApJ*, 588, 824
- Zhao, G., Zhao, Y., Chu, Y., Jing, Y., Deng, L., 2012, *RAA*, 12, 723